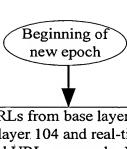Fig. 1

**Fig. 2**

Beginning of new epoch

302

Identify the base segment 112 of URLs from base layer 102 to be crawled during an epoch (*e.g.* on a given day). Refine daily layer 104 and real-time layer 106 based on how frequently URL content changes and URL page rank. Review history logs 218 to determine how frequently URL content changes. Refine URL crawl frequency and/or select URLs that are to be crawled.

304

Deliver the selected URLs to URL managers 204

Place real-time layer 106 in a separate URL manager 204 that does not scan link logs 214 or status logs 212

Scan link logs 214 to discover new URLs. Ignore URLs that have been seen before but are not currently scheduled for crawling. Add newly discovered URLs to the pool of URLs to be crawled

Scan status logs 212 to update the states of URLs that have been crawled during the given time period

306

Deliver URLs to URL server 206 when they are requested by server 206

Request reservations from host load servers that host URLs to be crawled

Distribute URLs to robots 208 in accordance with terms of the conditions for the reservations
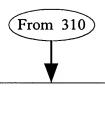
308

Crawl URLs

Perform local DNS resolution when possible using a local DNS resolution database

Provide support for cookies when needed using a cookie database

310

Eliminate duplicate URLs with Dupserver 224

To 312

Fig. 3A

From 310

*312

Filter data obtained from crawling

Write out links found in crawled pages to link logs 214

Write out document associated with each crawled URL indexed by page rank
to Rtlogs 226, 228 or 330 as a function of which layer the URL belongs to in
data structure 100

Write out one history record for each URL that was crawled to history logs 218,
where the history record includes information such as the URL, crawl status,
checksums, time taken to download, *etc*

Write out one status record for each URL that was crawled to status logs 212 that
contains information such as crawl status

*314

Real-time indexers 232 obtain documents from the Rtlogs and index them so
that they are searchable by a front-end querying system (not shown)

*316

Global state manager 216 creates (i) link maps 220 that do not contain text
506 and are keyed by the fingerprints of URLs 502 and (ii) anchor maps 238
that contain text 506 and are keyed by the fingerprints of outbound URLs 504

Fig. 3B

600

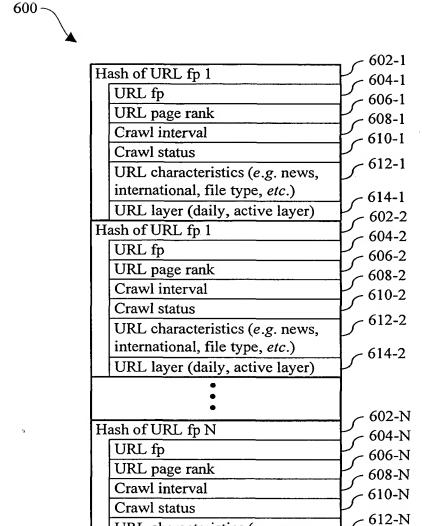| | |
|---|---|
| Hash of URL fp 1 | 602-1 |
| URL fp | 604-1 |
| URL page rank | 606-1 |
| Crawl interval | 608-1 |
| Crawl status | 610-1 |
| URL characteristics (*e.g.* news, international, file type, *etc.*) | 612-1 |
| URL layer (daily, active layer) | 614-1 |
| Hash of URL fp 1 | 602-2 |
| URL fp | 604-2 |
| URL page rank | 606-2 |
| Crawl interval | 608-2 |
| Crawl status | 610-2 |
| URL characteristics (*e.g.* news, international, file type, *etc.*) | 612-2 |
| URL layer (daily, active layer) | 614-2 |
| ⋮ | |
| Hash of URL fp N | 602-N |
| URL fp | 604-N |
| URL page rank | 606-N |
| Crawl interval | 608-N |
| Crawl status | 610-N |
| URL characteristics (*e.g.* news, international, file type, *etc.*) | 612-N |
| URL layer (daily, active layer) | 614-N |

Fig. 4

## Link Log

214

| Record for URL document 1 | 502-1 |
| --- | --- |
|    Outbound URL 1 | 504-1-1 |
|    Text around URL 1 | 506-1-1 |
|    • • • | |
|    Outbound URL X | 504-1-X |
|    Text around URL X | 506-1-X |
| • • • | |
| Record for URL document N | 502-N |
|    Outbound URL 1 | 504-N-1 |
|    Text around URL 1 | 506-N-1 |
|    • • • | |
|    Outbound URL Z | 504-N-Z |
|    Text around URL Z | 506-N-Z |

**Fig. 5A**

## Rtlog

226/228/230

| Pair 1 | 510-1 |
| --- | --- |
|    Document | 512-1 |
|    Page rank | 514-1 |
| Pair 2 | 510-2 |
|    Document | 512-2 |
|    Page rank | 514-2 |
| • • • | |
| Pair M | 510-3 |
|    Document | 512-3 |
|    Page rank | 514-3 |

**Fig. 5B**

## History Log

218

| Record for URL 1 | 520-1 |
| --- | --- |
|    URL fp | 522-1 |
|    Crawl status | 524-1 |
|    Content checksum | 526-1 |
|    Approx. content checksum | 528-1 |
|    Link checksum | 530-1 |
|    Source | 532-1 |
|    Time taken to download | 534-1 |
|    Error condition | 536-1 |
|    • • • | |
| Record for URL 2 | 520-2 |
| • • • | |

**Fig. 5C**

## Status Log

212

| Record for URL 1 | 550-1 |
| --- | --- |
|    URL | 522-1 |
|    URL fingerprint | 554-1 |
|    Crawl status | 524-1 |
|    Content checksum | 526-1 |
|    Outgoing links | 556-1 |
|      Outbound URL 1 | 504-1-1 |
|      Outbound URL 2 | 504-1-2 |
|    • • • | |
|    Duplicate status | 558-1 |
|    • • • | |
| Record for URL 2 | 550-2 |
| • • • | |

**Fig. 5D**

**Fig. 5**

Fig. 6

1012-1
1012-2
1012-Y

1010-1
1010-2
1010-Y

1002
1011

1008-1
1008-2
1008-X

1006

1000

1004-1
1004-2
1004-X

011378-0007-999

Fig. 7

011378-0007-999

Fig. 8

220

1110-M

1110-3
1110-2
1110-1

Sorted Link
Map

Merged Link
Map

1110-(M+1)

216

1103

Link Log

214

Link
Record
Sorter    1202

Anchor
Log      1206

Anchor
Sorter   1208

Link Map
Merger   1204

Anchor Map
Merger   1210

Merged
Anchor Map

1112-(N+1)

1112-3
1112-2
1112-1

Sorted Anchor
Map

1112-N

238

011378-0007-999

| | | |
|---|---|---|
| URL-1 | URL-1-1 | "what URL-1-1 says about URL-1" |
| | URL-1-2 | "what URL-1-2 says about URL-1" |
| | ⋮ | |
| | URL-1-N1 | "what URL-1-N1 says about URL-1" |
| URL-2 | URL-2-1 | "what URL-2-1 says about URL-2" |
| | URL-2-2 | DELETE URL-2-2 |
| | ⋮ | |
| | URL-2-N2 | "what URL-2-N2 says about URL-2" |
| ▪ ▪ ▪ | | |
| URL-M | URL-M-1 | "what URL-M-1 says about URL-M" |
| | URL-M-2 | "what URL-M-2 says about URL-M" |
| | ⋮ | |
| | URL-M-N1 | "what URL-M-N1 says about URL-M" |

1302-1
1302-2
1304-1
1304-2
1304-K
1302-M
1303
1112

Fig. 9

011378-0007-999

| URL-1 | URL-1-1 URL-1-2 ... URL-1-N1 | URL-2 | URL-2-1 URL-2-2 ... URL-2-N2 | ■ ■ ■ | URL-M | URL-M-1 URL-M-2 ... URL-M-N1 |

1402-1

1402-2

1404-1
1404-2

1404-N2

1402-M

1403

1110

Fig. 10

Fig. 11

Fig. 12